

# MusicWeb: music discovery with open linked semantic metadata

Mariano Mora-Mcginity<sup>1</sup>, Alo Allik<sup>1</sup>, György Fazekas<sup>1</sup>, and Mark Sandler<sup>1</sup>

Queen Mary University, London, UK

{m.mora-mcginity, a.allik, g.fazekas, mark.sandler}@qmul.ac.uk

**Abstract.** This paper presents MusicWeb, a novel platform for music discovery by linking music artists within a web-based application. MusicWeb provides a browsing experience using connections that are either extra-musical or tangential to music, such as the artists’ political affiliation or social influence, or intra-musical, such as the artists’ main instrument or most favoured musical key. The platform integrates open linked semantic metadata from various Semantic Web, music recommendation and social media data sources. Artists are linked by various commonalities such as style, geographical location, instrumentation, record label as well as more obscure categories, for instance, artists who have received the same award, have shared the same fate, or belonged to the same organisation. These connections are further enhanced by thematic analysis of journal articles, blog posts and content-based similarity measures focussing on high level musical categories.

**Keywords:** Semantic Web, Linked Open Data, music metadata, semantic audio analysis, music information retrieval

## 1 Introduction

In recent years we have witnessed an explosion of information, a consequence of millions of users producing and consuming web resources. Researchers and industry have recognised the potential of this data, and have endeavoured to develop methods to handle such a vast amount of information. There are two main approaches to music recommendation [1]: the first is known as *collaborative filtering* [2], which recommends music items based on the choices of similar users. The second model is based on audio content analysis, or *music information retrieval*. The task here is to extract low to high-level audio features such as tempo, key, metric structure, melodic and harmonic sequences, instrument recognition and song segmentation, which are then used to measure music similarity. There are, however, limitations in both approaches to music recommendation. Most users participating in collaborative filtering listen to a very small percentage of the music available, the so called “short-tail”, whereas the much larger “long-tail” remains mainly unknown [3]. Many music listeners follow artists because of their style and would be interested in music from similar artists. There are many different ways in which people are attracted to new artists: word of mouth, their network of friends, music magazines or blogs, songs heard in a

movie or a T.V. commercial, they might be interested in a musician who has played with another artist or been mentioned as an influence, etc. The route from listening to one artist and discovering a new one would sometimes seem very disconcerting were it to be drawn on paper. A listener is not so much following a map as exploring new territory, with many possible forks and shortcuts. Music discovery systems generally disregard this kind of information, often because it is very nuanced and difficult to parse and interpret.

MusicWeb is a music discovery platform which offers users the possibility of exploring editorial, cultural and musical links between artists. It gathers, extracts and manages metadata from many different sources, including DBpedia.org, Sameas.org, MusicBrainz, the Music Ontology, Last.FM and Youtube as well as editorial and content-derived information. The connections between artists are based on YAGO categories [4], which successfully extracts categories from each wikipedia entry after contrasting them with WordNet. These are categories such as style, geographical location, instrumentation, record label, but also more obscure links, for instance, artists who have received the same award, have shared the same fate, or belonged to the same organisation or religion. These connections are further enhanced by thematic analysis of journal articles, blog posts and content-based similarity measures focusing on high level musical categories.

## 2 MusicWeb Architecture

MusicWeb provides a browsing experience using connections that are either extra-musical or tangential to music, such as the artists' political affiliation or social influence, or intra-musical, such as the artists' main instrument or most favoured musical keys. It does this by pulling data from several different web knowledge content resources and presenting them for the user to navigate in a faceted manner [5]. The listener can begin his journey by choosing or searching an artist. The application offers Youtube videos, audio streams, photographs and album covers, as well as the artist's biography (see example in Fig. 1). The page also includes many box widgets with links to artists who are related to the current artist in different, and sometimes unexpected ways.

MusicWeb was originally conceived as a platform for collating metadata about music artists using already available online linked data resources. The core functionality of the platform relies on available SPARQL endpoints as well as various commercial and community-run application programming interfaces (APIs).

The MusicWeb API uses a number of linked open data (LOD) resources and Semantic Web ontologies to process and aggregate information about artists:

**Musicbrainz**<sup>1</sup> is an online, open, crowd-sourced music encyclopedia, that provides reliable and unambiguous identifiers for entities in music publishing metadata, including artists, releases, recordings, performances, etc.

**DBPedia**<sup>2</sup> is a crowd-sourced community effort to extract structured informa-

---

<sup>1</sup> <http://musicbrainz.org>

<sup>2</sup> <http://dbpedia.org>



Fig. 1. Example of a MusicWeb artist page.

tion from Wikipedia and make it available on the Web.

**Sameas.org**<sup>3</sup> manages Universal Resource Identifier (URI) co-references on Web of Data.

**Youtube** API is used to query associated video content for the artist panel.

**Last.fm**<sup>4</sup> is an online music social network and recommender system that collects information about users listening habits and makes available crowd-sourced tagging data through an API.

**YAGO** is a semantic knowledge base that collates information and structure from Wikipedia, WordNet and GeoNames.

**The Music Ontology** [6] provides main concepts and properties for describing musical entities, including artists, albums, tracks, performances, compositions, etc., on the Semantic Web

The global MusicBrainz identifiers enable convenient and concise means to disambiguate between potential duplicates or irregularities in metadata across resources, a problem which is all too common in systems relying on named entities. Besides identifiers, the MusicBrainz infrastructure is also used for the search functionality of MusicWeb. However, in order to query any information in DBpedia, the MusicBrainz identifiers need to be associated with a DBpedia resource, which is a different kind of identifier. This mapping is achieved by querying the Sameas.org co-reference service to retrieve the corresponding DBpedia URIs.

<sup>3</sup> <http://sameas.org>

<sup>4</sup> <http://last.fm>

The caveat in this process is that Sameas does not actually keep track of MusicBrainz artist URIs, however, by substituting the domain for the same artist's URI in the BBC domain<sup>5</sup>, MusicWeb can get around this obstacle. Once the DBpedia artist identity is determined, the service proceeds to construct the majority of the profile, including the biography and most of the linking categories to other artists. The standard categories available include associated artists and artists from the same hometown, while music group membership and artist collaboration links are queried from MusicBrainz. The core of the Semantic Web linking functionality is provided by categories from YAGO. The Last.fm API provides with information on artists it deems similar.

### 3 Artist Similarity

There are many ways in which artists can be considered related: similarity may be based on a particular style or genre, but it may also mean that artists are followed by people from similar social backgrounds, political inclinations, or age groups. Artists can also be associated because they have collaborated, participated in the same event, or their lyrics touch upon similar themes. Linked data facilitates faceted searching and displaying of information [7]: an artist may be similar to many other artists in one of the ways just mentioned, and to a completely different plethora of artists in other senses, all of which might contribute to music discovery. Semantic Web technologies can help us gather different facets of data and shape them into representations of knowledge. MusicWeb does this by searching similarities in three different domains: socio-cultural, research and journalistic literature and content-based information retrieval.

**Socio-cultural** connections between artists in MusicWeb are primarily derived from YAGO categories that are incorporated into entities in DBpedia. Many categories, in particular those that can be considered extra-musical or tangential to music, stem from the particular methodology used to derive YAGO information from Wikipedia. While DBpedia extracts knowledge from the same source, YAGO leverages Wikipedia category pages to link entities without adapting the Wikipedia taxonomy of these categories [4]. The hierarchy is created by adapting the Wikipedia categories to the WordNet concept structure. This enables linking each artist to other similar artists by various commonalities such as style, geographical location, instrumentation, record label as well as more obscure categories, for example, artists who have received the same award, have shared the same fate, or belonged to the same organisation or religion. YAGO categories can reveal connections between artists that traditional isolated music datasets would not be able to establish.

**Literature-based** linking is achieved by data-mining research articles and online publications using natural language processing. MusicWeb uses Mendeley<sup>6</sup> and Elsevier<sup>7</sup> databases for accessing research articles that are curated and cate-

---

<sup>5</sup> <http://www.bbc.co.uk/music/artists/>

<sup>6</sup> <http://dev.mendeley.com/>

<sup>7</sup> <http://dev.elsevier.com/>

gorised by keywords, authors and disciplines. Online publications, such as newspapers, music magazines and blogs focused on music, on the other hand, constitute non-curated data. Relevant information in this case must be extracted from the body of the text. The data is collated by crawling websites by keywords or tags in the title and by following external links contained in pages. Many texts contain references to an artist name without actually being relevant to MusicWeb. A search for Madonna, for example, can yield many results from the fields of sculpture, art history or religion studies. The first step is to model the relevance of the text, and discard texts which are of no interest to music discovery. Texts and abstracts are then subjected to semantic analysis. The text as a bag of words is used to query the Alchemy<sup>8</sup> language analysis service for named entity recognition and keyword extraction. The entity recogniser provides a list of names that appear mentioned in the text together with a measure of relevance. MusicWeb identifies musical artists using its internal artist database as well as DBpedia, MusicBrainz and Freebase. Keyword extraction is used for non-curated sources and involves checking keywords against WordNet for hypernyms. Artists that share keywords or hypernyms are considered to be relevant to the same topic in the literature. MusicWeb also offers links between artists who appear in different articles by the same author, as well as in the same journal.

**Content-based linking** involves methodology of Music Information Retrieval (MIR) [8] which facilitate applications that rely on perceptual, statistical, semantic or musical features derived from audio using digital signal processing and machine learning methods. These features may include statistical aggregates computed from time-frequency representations extracted over short time windows. Higher-level musical features include keys, chords, tempo, rhythm, as well as semantic features like genre or mood, with specific algorithms to extract this information from audio. High-level stylistic descriptors can correlate with lower level features such as the average tempo of a track, the frequency of note onsets, the most commonly occurring keys or chords or the overall spectral envelope that characterises instrumentation. To exploit different types of similarity, we model each artist using three main categories of audio descriptors: rhythmic, harmonic and timbral. We compute the joint distribution of several low-level features in each category over a large collection of tracks from each artist. We then link artists exhibiting similar distributions of these features. The features are obtained from the AcousticBrainz<sup>9</sup> Web service which provides descriptors in each category of interest. Tracks are indexed by MusicBrainz identifiers enabling unambiguous linking to artists and other relevant metadata. For each artist in our database, we retrieve features for a large collection of their tracks in the above categories, including beats-per-minute and onset rate (rhythmic), chord histograms (harmonic) and Mel-Frequency Cepstral Coefficients (timbral) features.

---

<sup>8</sup> AlchemyAPI is used under license from IBM Watson.

<sup>9</sup> <https://acousticbrainz.org/>

## 4 Conclusions

MusicWeb is an emerging application to explore the possibilities of linked data-based music discovery. The methods of linking artists employed in the system are intended to overcome issues such as infrequent access of lesser known artists in large music catalogues (the “long tail” problem) or the difficulty of recommending artists without user ratings in systems that employ collaborative filtering (“cold start” problem) [3]. This facilitates users to engage in interesting discovery paths through the space of music artists. Although similar to recommendation, this is in contrast with most recommender systems which operate on the level of individual music items. We aim at creating links between artists based on stylistic elements of their music derived from a collection of recordings and complement the social and cultural links. Future work will address investigating various different approaches to music discovery and how they can benefit from linked music metadata. The next steps are directed towards evaluating the potential acceptance of MusicWeb by end users to find out which linking methods listeners find appealing or interesting, and which they would use most.

## References

- [1] Yading Song, Simon Dixon, and Marcus Pearce. A survey of music recommendation systems and future perspectives. In *9th International Symposium on Computer Music Modeling and Retrieval*, 2012.
- [2] S Sneha, DS Jayalakshmi, J Shruthi, and Uttarika Ratnakar Shetty. Recommending music by combining content-based and collaborative filtering with user preferences. In *Emerging Research in Electronics, Computer Science and Technology*, pages 507–515. Springer, 2014.
- [3] Ò. Celma. *Music Recommendation and Discovery: The Long Tail, Long Fail, and Long Play in the Digital Music Space*. Springer Verlag, 2010.
- [4] MS Fabian, K Gjergji, and W Gerhard. Yago: A core of semantic knowledge unifying wordnet and wikipedia. In *16th International World Wide Web Conference, WWW*, pages 697–706, 2007.
- [5] Gary Marchionini. Exploratory search: from finding to understanding. *COMMUNICATIONS OF THE ACM*, 49(9), 2006.
- [6] Yves Raimond, Samer A Abdallah, Mark B Sandler, and Frederick Giasson. The music ontology. In *ISMIR*, pages 417–422. Citeseer, 2007.
- [7] Miguel Ángel Rodríguez-García, Luis Omar Colombo-Mendoza, Rafael Valencia-García, Antonio A Lopez-Lorca, and Ghassan Beydoun. Ontology-based music recommender system. In *Distributed Computing and Artificial Intelligence, 12th International Conference*, pages 39–46. Springer, 2015.
- [8] Michael A. Casey, Remco Veltkamp, Masataka Goto, Marc Leman, Christophe Rhodes, and Malcolm Slaney. Content-based music information retrieval: Current directions and future challenges. volume 96. IEEE Proceedings, April 2008.